

Actuarial Data Science: Opportunities and Challenges

Junge Aktuare, Swiss Association of Actuaries

8th March 2022

Dr. Jürg Schelldorfer, Actuary SAA

Senior Actuarial Data Scientist, Swiss Re

Chair of the «Data Science» working group of the Swiss Association of Actuaries (SAA)

Disclaimer

The opinions expressed in this presentation are those of the author only. They are inspired by the work that the author is doing for both Swiss Re and the SAA, but they do not necessarily reflect any official view of either Swiss Re or the SAA.

Table of Content

1. Introduction
2. Opportunities
3. Challenges
4. Take-home message

Introduction

Major Topics

FINMA-regulated insurance market



IFRS 17

Actuarial Data
Science

Regulation

Low interest
rates

Sustainability

SAV Fachgruppe «Data Science»

- Anja Friedrich
 - Frank Genheimer
 - Thomas Hull
 - Dr. Christian Lorentzen
 - David Lüthi (Stv)
 - Dr. Michael Mayer
 - Dr. Daniel Meier (Stv)
 - Dr. Jürg Schelldorfer (Leitung)
 - Dr. Alessandro Torre
 - Dr. Andreas Troxler
 - Prof. Dr. Mario Wüthrich
- ...und viele weitere welche in den letzten 5 Jahren temporär mitgearbeitet haben.

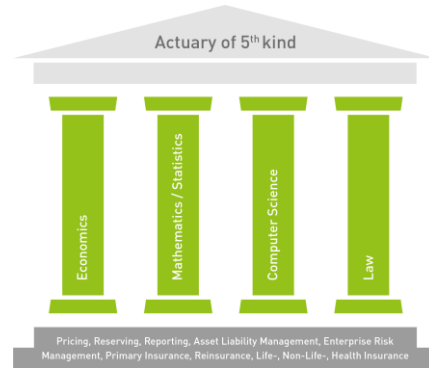
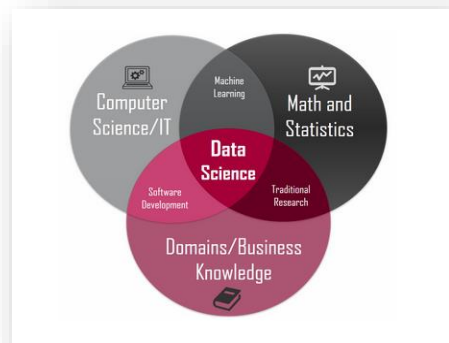
www.actuarialdatascience.org



A screenshot of the website homepage. At the top left is the logo for SAV (Schweizerische Aktuarvereinigung), ASA (Association Suisse des Actuaires), and ASA (Associazione Svizzera degli Attuari). To the right of the logo is the text 'Actuarial Data Science' and 'An Initiative of the Swiss Association of Actuaries'. Below this is a navigation menu with links: Home, ADS Tutorials, ADS Strategy, ADS Lectures / Courses, ADS Regulatory / Ethics, DS Lectures / Books, External Courses, Newsletter, and About Us. The main content area has a 'Home' section with the text: 'The main purpose of this website is to make the work and results of the working group "Data Science" of the Swiss Association of Actuaries (SAA) / Schweizerische Aktuarvereinigung (SÄV) easily available to interested people. Actuarial Data Science (ADS) is defined to be the intersection of Actuarial Science (AS) and Data Science (DS). The core targets are: • ADS Tutorials: Writing tutorials for actuaries which provide a thorough and yet easy introduction to various methods from Data Science. We provide methodological papers together with the code, such that everyone can easily learn the methods on his own data. • ADS Strategy: We have worked out a strategy for the Swiss Association of Actuaries.' To the right is an 'Updates' section with the text: 'Below, we provide the most recent changes to the website: • 19th July 20: Publication of our ninth tutorial: Convolutional neural network studies: (1) anomalies in mortality rates (2) image recognition (incl. code) • 7th May: Publication of our eighth tutorial: Peeking into the Black Box: An'.

What is (Actuarial) Data Science?

Definition(s) und differences Data Science / Actuarial Science⁽¹⁾



	Actuarial Science	Data Science
Basics	Mathematical Basics	
Data	Small Data	Small and Big Data
	Structured & static Data	Unstructured & dynamic Data
	Internal Data	External Data
Mathematics & Statistics	Probability Theory	Computational Statistics
	Life and Non-Life Insurance Mathematics	Algorithm
	Quantitative Risk Management	Information Theory
Computer science		Machine Learning & Visualisations
		Numeric Optimization
		Data Management
Programming languages	SAS, S Plus, R	Python, R
	SQL	SQL
	Excel / VBA	Julia, Spark, Scala
Domains/ business knowledge	Reserving, Pricing	
	ERM, ALM, Solvency	
	Accounting, Economics, Law	

Typical course content and competences



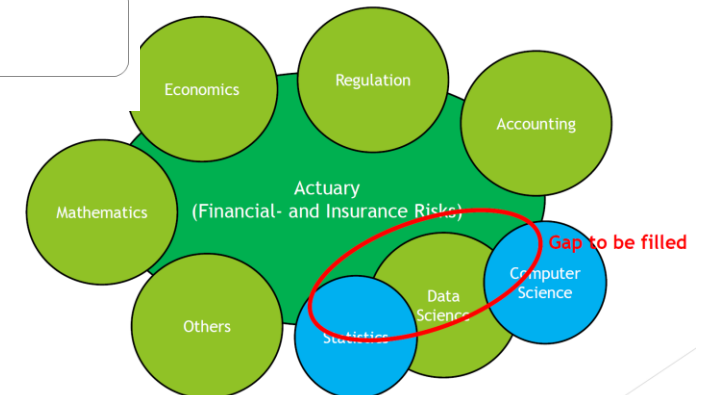
DATA SCIENCE STRATEGY

Data Science working group of the Swiss Association of Actuaries (SAA)

Version 2.0, August 2018
Published version

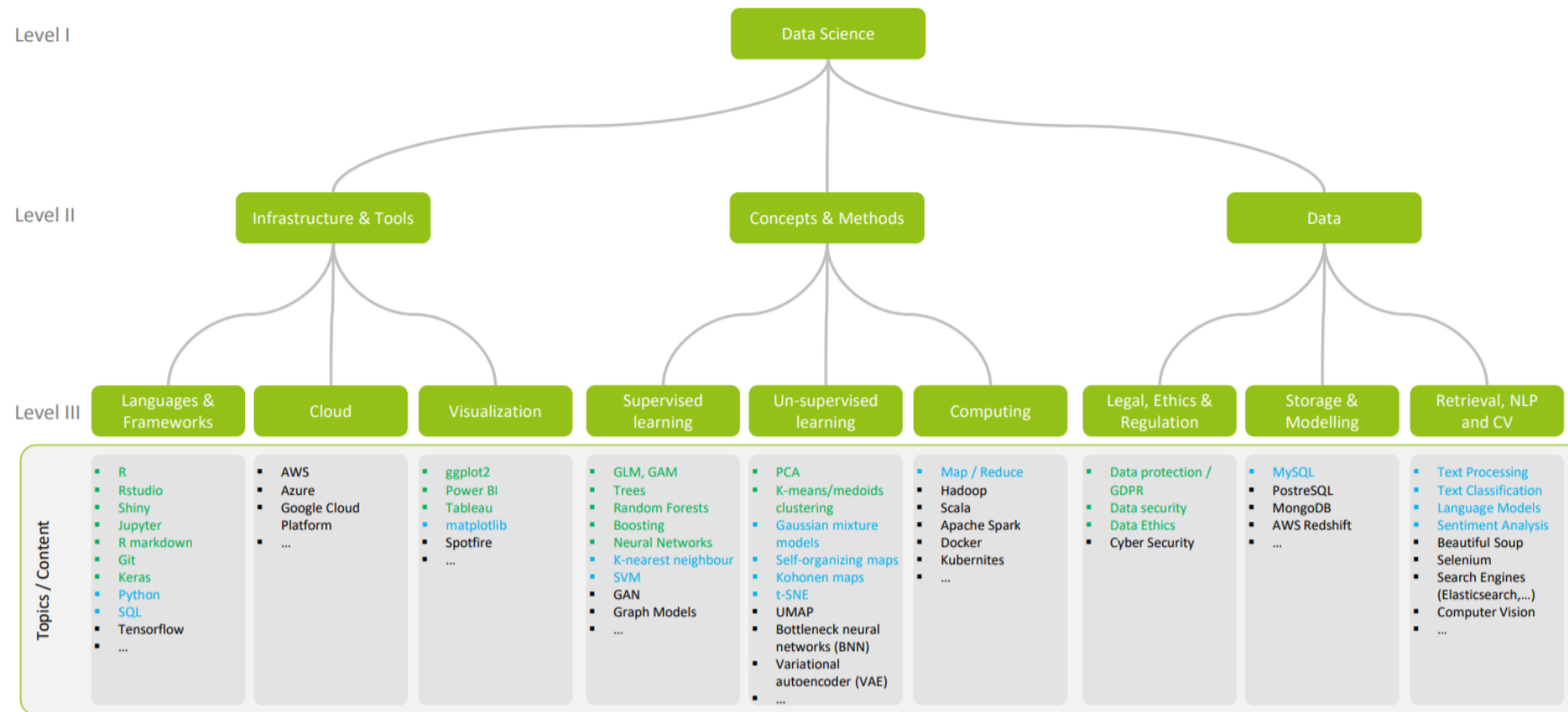
... intersection of ...

- Mathematics and statistics
- Computer science / IT
- Domain / Business Knowledge



⁽¹⁾ Aus: [Data Science Strategy, Data Science working group, SAA, Version 2.0, August 2018](#)

What is Data Science?



Green: Recommended for actuaries in the industry with some basic knowledge in data science

Blue: Recommended for actuaries in the industry with extended knowledge in data science (combined with green)

Black: Data science (combined with green and blue)

Opportunities

1 – Unstructured data

- Actuaries are trained and used to work with tabular/structured data. The big(ger) amount of data in an insurance company are unstructured data (text, images, pdfs,...)

By **structured data**, we mean data⁽⁷⁾ "organised into a formatted repository, typically a [relational] database", where it is stored in rows and columns. **Unstructured data**⁽³²⁾ is digitised information that is not organised in a pre-defined format. Examples include image, text, video and voice data.

- Unstructured data should and can be used by actuaries, as many technologies (e.g. extract text from pdf's) become a «commodity».
- **Concrete Examples:**
 - Free text fields to classify a claims to a specific claim type (e.g. glass, theft,...)
 - Extract information from structured pdf's
 - Extract daily COVID-19 cases from pdf's
 - Use text field to model the claims severities for worker's compensation claims data (Tutorial [here](#))
 - Predict number of injured in car collisions from police reports <-> insurance claims (Tutorial to come)

2 – More data

- Having big data is not the standard for actuaries, it is still the exception. Insurance companies are not the big data owners, as they are in the second line of the industry (enabling business, not creating business).
- Design products and actively engage in product development. Ideas: [The Geneva Association](#), [Swiss Re](#)
- Collaborations with third-party data provider
- **Concrete Examples:**
 - Telematics data
 - Vessel real-time data ([Link](#))
 - Shipment goods data
 - Sensor data for property insurance ([The Geneva Association: From Risk Transfer to Risk Mitigation](#))
 - Electronic Health Records (EHR) data
 - (Tracking data)
 - (...)

3 – Models and algorithms

- Advanced statistical and machine learning algorithms are **seamlessly** available through open source software libraries (e.g. Python, R)
- Commercial software providers have recognized the first point.
- Neural networks can be easily fitted with Python/R using a few lines of code.
- New is the seamlessly availability and the simple usage of the algorithms. The statistical models are not new.

Concrete Examples:

- Fitting a simple Neural Network using R/Python
- Fitting Classification and Regression Tree, Random Forest, Boosting
- Fitting clustering algorithms

R interface to Keras



CRAN Task View: Machine Learning & Statistical Learning
Maintainer: Torsten Hothorn
Contact: Torsten.Hothorn@r-project.org
Version: 2022-03-02
URL: <https://cran.r-project.org/view=MachineLearning>

Several add-on packages implement ideas and methods developed at the borderline between computer science and statistics:

- **Neural Networks and Deep Learning**: Single-hidden-layer neural networks are implemented in package `neuralnet`, (denoising autoencoder, stacked denoising autoencoder, restricted Boltzmann machine, see `rbm` and `M5`). The `Cubist` package fits rule-based models (similar to trees) with linear regression models in two recursive partitioning algorithms with unbiased variable selection and statistical stopping criteria. Distributions of the response are available in package `party` and `partykit` as well. Graphical tools for the visualization of trees are available in package `party`. Partitioning of mixture models is performed by `RPM`.
- **Computational infrastructure for representing trees and unified methods for prediction and visualization**: The reference implementation of the random forest algorithm for regression and c -inference trees is implemented in package `party`. `randomForestSRC` implements a unified treatment of random forest algorithms. In addition, packages `mapet` and `Rborist` offer R interfaces to fast C++ packages `RG` is an interface to a Python implementation of a procedure called regularized greedy fit.
- **Regularized and Sparse Models**: Regression models with some constraint on the parameter estimates (generalized linear models and Cox models) can be obtained from functions available in package `glmnet`. Package `Rx` can be used to identify and display TRACs for a specified shrinkage p -optimal LASSO penalty to produce sparse solutions is implemented in package `penalized`. The interface to the LIBLINEAR library. The `lasso2` package fits linear and logistic regression model errors is estimated by `lasso`, inference on low-dimensional components of Lasso regression and of c are fitted by package `penalized` using composite optimization by conjugation operator. The `lasso2` package `Boosting and Gradient Descent`: Various forms of gradient boosting are implemented in package `gb`.
- **Generalized linear, additive and nonparametric models** is available in package `uboot`.
- **Support Vector Machines and Kernel Methods**: The function `svm()` from `e1071` offers an interface to spaces can be estimated using `libsvm`, which also offers procedures for model selection and prediction.
- **Bayesian Methods**: Bayesian Additive Regression Trees (BART), where the final model is defined by package `bart`, Bayesian structure learning in undirected graphical models for multivariate continue `Optimization using Genetic Algorithms: Package rgenoud offers optimization routines based on genetic algorithms.`
- **Association Rules**: Package `arules` provides both data structures for efficient handling of sparse binary data efficiently, in the form of self-sufficient streams, using either leverage or lift.
- **Fuzzy Rule-based Systems**: Package `fuzzy` implements a host of standard methods for learning fuzzy rules.
- **Model selection and validation**: Package `e1071` has function `train()` for hyper parameter tuning and other visualization techniques for comparing candidate classifiers are available from package `Rkde`.
- **Causal Machine Learning**: The package `DoubleML` is an object-oriented implementation of the don't
- **Other procedures**: Evidential classifiers quantify the uncertainty about the class of a test pattern using `Meta packages`: Package `caret` provides miscellaneous functions for building predictive models, including similar toolboxes. The `mlr` package implements a general purpose machine learning platform that has machine learning algorithms, such as nearest neighbors, trees, random forests, and several feature sets.
- **GUI** `mlr` is a graphical user interface for data mining in R.
- **Visualization (initially contributed by Brandon Greenwell)**: The `stats::tsnePlot()` function package rendering of the prediction function, are implemented in a few packages: `glm`, `randomForest` and `g` PDPs for a wide variety of machine learning models (e.g., random forests, support vector machines, less information. `K-Flow` focuses on constructing individual conditional expectation (ICE) curves, a

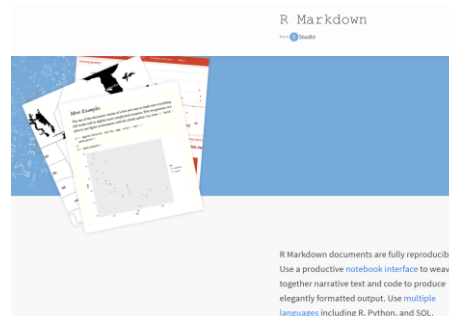
<https://keras.rstudio.com/>

4 – Technology and computing power

- «2 – Big(ger) data» and «3 - Algorithms» are useful due to the available computing power to fit complex mathematical models.
- Many technologies/tools (see below) enable the development of machine learning models

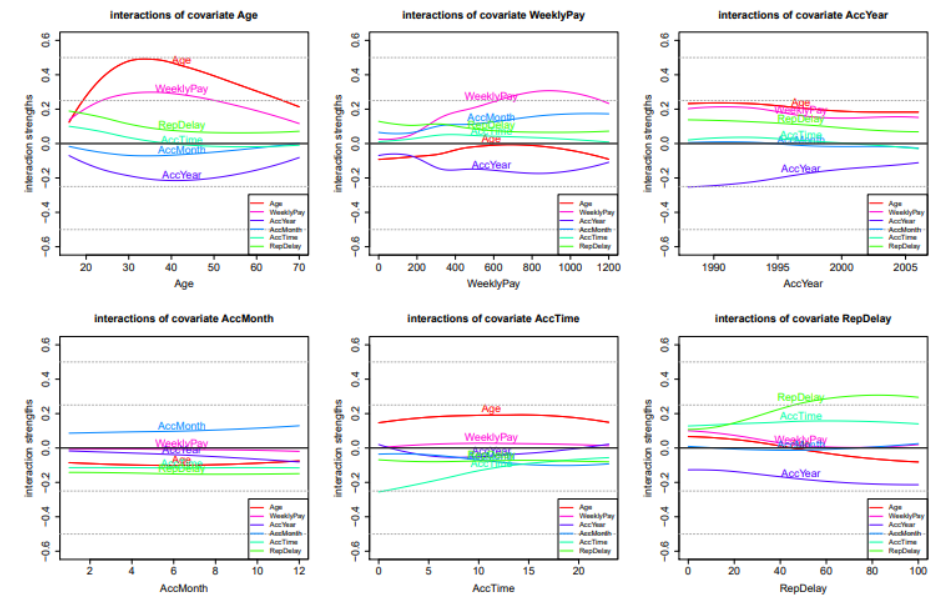
- **Concrete Examples:**

- Git ([Link](#))
- Jupyter notebooks ([Link](#))
- R Shiny Apps ([Link](#)): Deploy a dashboard company-internally for management
- R Markdown ([Link](#))
- Deploy and retrain models frequently
- Deploy pricing model, re-trained every day ([Parametric Flight Delay Insurance](#), Swiss Re)
- Accident images
- Satellite images



5 – Processes and efficiency

- Using new technological tools (see previous slides) enables actuaries to improve processes and increase their efficiency, and allows the actuaries to focus on the core.
- Concrete Examples:
 - R package(s) for calculating special reserves
 - Automating reports with R markdown ([EAA training](#))
 - Find easily structure not yet caputed in a simple model (Tutorial [here](#))
 - Speed up finding interactions in glms



Challenges

1 – Data quality

- Data are the foundation of data science.
- The quality of the data remains as important as ever! «Garbage in – garbage out». If the data are not meeting the quality criteria (accuracy, appropriateness, completeness), there is no reason to draw any business conclusions from it.
- See Solvency II data quality definition.
- Data Governance
- Collecting / ensuring data quality can not be easily achieved → Data strategy
- Statistical / machine learning techniques can help to detect anomalies in data (anomaly, outlier detection).
- **Concrete Examples:**
 - Predict data points and compare with observed value (missForest, [missRanger](#))
 - Data strategy
 - Data governance

mayer79/ missRanger

R package "missRanger" for fast imputation of missing values by random forests.

4 Contributors 2 Issues 41 Stars 10 Forks



2 – Usage and access to data

- Usage of some data is not allowed (legal, regulation)
- Having the appropriate / desired data is not the case
- Collaborations with other companies
- Buying external data

- **Concrete Examples:**
 - Submission data for a corporate insurer / reinsurer
 - The insured is the data owner, or has the granular data → information asymmetry

3 – XFT

- Product development means working in cross-functional teams (XFT), where actuaries / data specialists need to help design and define a product.
- Just doing the math/pricing is not enough.
- Interest beyond the actuarial core is required.
- Data quality and data collection needs time and care.

- **Concrete Examples:**
 - Shipment goods data
 - Swiss Re flight delay insurance
 - Swiss Re P&C Analytics ([Link](#)), see Impact+

4 – Varia

- Legacy in IT infrastructure and systems
- Legacy in mindset and lack of knowledge about opportunities/limits of Actuarial Data Science
- Sharing/Publish data is difficult than 10 years ago
- Sharing use cases is often not of interest
- Better/more complex algorithms are rarely the biggest challenge for an innovation.

- **Concrete Examples:**
 - No public text data on insurance claims
 - For learning, synthetic data are sufficient.
 - Synthetic data generators ([Simulation Machine](#), ...)

Take-home message

Conclusions

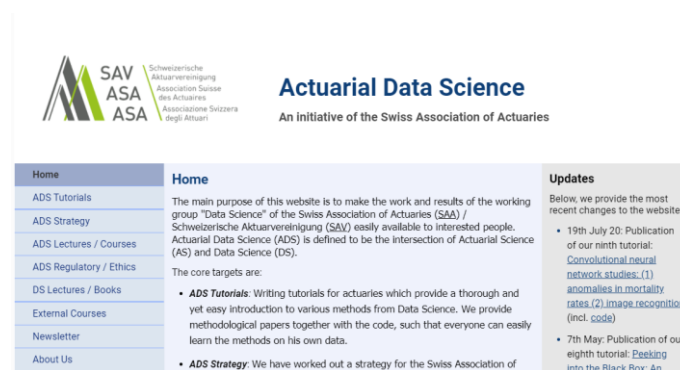
- Data Science != Actuarial Science
- Actuaries will collaborate with data scientists as they collaborate with IT, accounting, claims, underwriting,...
- Opportunities for additional data, better models and better toolkits
- Data quality/usage/access is a major challenge
- A very well calibrated GLM may still be as good as an advanced machine learning model in terms of accuracy.

Visit

For actuaries and data scientists in insurance



www.actuarialdatascience.org



Actuarial Data Science
An initiative of the Swiss Association of Actuaries

Home

- ADS Tutorials
- ADS Strategy
- ADS Lectures / Courses
- ADS Regulatory / Ethics
- DS Lectures / Books
- External Courses
- Newsletter
- About Us

Updates

Below, we provide the most recent changes to the website:

- 19th July 20: Publication of our ninth tutorial: [Convolutional neural network studies. \(1\) anomalies in mortality rates. \(2\) image recognition \(incl. code\)](#)
- 7th May: Publication of our eighth tutorial: [Peeking into the Black Box: An](#)



Actuarial Data Science in the Swiss Association of Actuaries
Insurance · Zurich, Zurich · 1,671 followers

An initiative of the Swiss Association of Actuaries and its Data Science Working Group.

[Visit website](#)

Acknowledgements

People:

- [All members of the SAA working group](#)
- Dr. Alexander Noll
- Dr. Simon Renzmann
- Ron Richman

Institutions:

- [Swiss Association of Actuaries \(SAA\)](#)
- [RiskLab at ETH Zurich](#)
- [MobiLab for Analytics at ETH Zurich](#)

Companies:

- [Swiss Re](#)

Appendix

ADS basics: Articles and repositories

The following articles/repositories are fundamental for entering the topic of Actuarial Data Science (ADS):

- [Data Analytics for Non-Life Insurance Pricing](#), ETH Zurich, [M.V. Wüthrich](#) and C. Buser
- [AI in Actuarial Science](#), R. Richman, SSRN, 2018
- [ADS Tutorials](#), SAA, 2018-present
- [Insurance Analytics – A Primer](#), International Summer School of the Swiss Association of Actuaries, 2018
- [Insurance Data Science: Use and Value of Unusual Data](#), International Summer School of the Swiss Association of Actuaries, 2019

And do not forget the fundamentals of Statistics vs. Machine Learning:

- [Statistical Modeling: The Two Cultures](#). L. Breimann, Statistical Science 16/3, 199-215, 2001
- [To explain or to Predict?](#), G. Shmueli, Statistical Science 25/3, 289-310, 2010